

5 Mutual Information and Channel Capacity

In Section 2, we have seen the use of a quantity called entropy to measure the amount of randomness in a random variable. In this section, we introduce several more information-theoretic quantities. These quantities are important in the study of Shannon's results.

5.1 Information-Theoretic Quantities

Definition 5.1. Recall that, the **entropy** of a **discrete** random variable X is defined in Definition 2.41 to be

$$H(X) = - \sum_{x \in S_X} p_X(x) \log_2 p_X(x) = -\mathbb{E}[\log_2 p_X(X)]. \quad (16)$$

Similarly, the entropy of a discrete random variable Y is given by

$$H(Y) = - \sum_{y \in S_Y} p_Y(y) \log_2 p_Y(y) = -\mathbb{E}[\log_2 p_Y(Y)]. \quad (17)$$

In our context, the X and Y are input and output of a discrete memoryless channel, respectively. In such situation, we have introduced some new notations in Section 3.1:

$$\begin{array}{lll} S_X \equiv \mathcal{X} & p_X(x) \equiv p(x) & p_{Y|X}(y|x) \equiv Q(y|x) \longrightarrow Q \text{ matrix} \\ S_Y \equiv \mathcal{Y} & p_Y(y) \equiv q(y) & p_{X,Y}(x,y) \equiv p(x,y) \longrightarrow p \text{ matrix} \end{array}$$

Under such notations, (16) and (17) become

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x) = -\mathbb{E}[\log_2 p(X)] \equiv H(\underline{p}) \quad (18)$$

the pmf of X expressed using a row vector.

and

$$H(Y) = - \sum_{y \in \mathcal{Y}} q(y) \log_2 q(y) = -\mathbb{E}[\log_2 q(Y)]. \quad (19)$$

Definition 5.2. The **joint entropy** for two random variables X and Y is given by

$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log_2 p(x, y) = -\mathbb{E}[\log_2 p(X, Y)].$$

$\underbrace{\sum_{(x,y)}}_{\text{elements inside the } p \text{ matrix}}$
 45

Example 5.3. Random variables X and Y have the following joint pmf matrix \mathbf{P} :

$$\mathbf{P} = \begin{array}{c} \begin{array}{c} \cancel{x \setminus y} \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} \begin{array}{c} y_1 \\ y_2 \\ y_3 \\ y_4 \end{array} \end{array} \begin{bmatrix} \frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{4} \\ \frac{1}{16} & \frac{1}{8} & \frac{1}{16} & 0 \\ \frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \\ \frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \end{bmatrix} \begin{array}{c} p(x) \\ \xrightarrow{\Sigma} \frac{1}{2} \\ \xrightarrow{\Sigma} \frac{1}{4} \\ \xrightarrow{\Sigma} \frac{1}{4} \\ \xrightarrow{\Sigma} \frac{1}{8} \end{array}$$

$\begin{array}{c} \Sigma \downarrow \\ \frac{1}{4} \end{array} \quad \begin{array}{c} \Sigma \downarrow \\ \frac{1}{4} \end{array} \quad \begin{array}{c} \Sigma \downarrow \\ \frac{1}{4} \end{array} \quad \begin{array}{c} \Sigma \downarrow \\ \frac{1}{4} \end{array}$

Find $H(X)$, $H(Y)$ and $H(X, Y)$.

$$H(X) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - 2 \times \frac{1}{8} \log_2 \frac{1}{8} = \frac{1}{2} + \frac{2}{4} + \frac{6}{8} = \frac{7}{4} \text{ bits}$$

$$H(Y) = -4 \times \frac{1}{4} \log_2 \frac{1}{4} = 2 = \log_2 4 = 2$$

$$H(X, Y) = -2 \times \frac{1}{8} \log_2 \frac{1}{8} - 6 \times \frac{1}{16} \log_2 \frac{1}{16} - \frac{1}{4} \log_2 \frac{1}{4} - 4 \times \frac{1}{32} \log_2 \frac{1}{32} - 3 \times 0 \log_2 0 = \frac{27}{8} \text{ bits}$$

Definition 5.4. Conditional entropy:

- (a) The (conditional) entropy of Y when we know $X = x$ is denoted by $H(Y|X = x)$ or simply $H(Y|x)$. It can be calculated from

$$H(Y|x) = - \sum_{y \in \mathcal{Y}} Q(y|x) \log_2 Q(y|x) = -\mathbb{E}[\log_2(Q(Y|x)) | X = x].$$

$\xrightarrow{\text{updated (revised)}}$

- Note that the above formula is what we should expect it to be. When we want to find the entropy of Y , we use (19):

$$H(Y) = - \sum_{y \in \mathcal{Y}} q(y) \log_2 q(y).$$

When we have an extra piece of information that $X = x$, we should update the probability about Y from the unconditional probability $q(y)$ to the conditional probability $Q(y|x)$.

- Note that when we consider $Q(y|x)$ with the value of x fixed and the value of y varied, we simply get the whole x -row from \mathbf{Q} matrix. So, to find $H(Y|x)$, we simply find the “usual” entropy from the probability values in the row corresponding to x in the \mathbf{Q} matrix.
- (b) The (average) conditional entropy of Y when we know X is denoted by $H(Y|X)$. It can be calculated from

$$\begin{aligned}
 H(Y|X) &= \sum_{x \in \mathcal{X}} p(x) H(Y|x) \\
 &= - \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} Q(y|x) \log_2 Q(y|x) \\
 &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log_2 Q(y|x) \\
 &= -\mathbb{E}[\log_2 Q(Y|X)]
 \end{aligned}$$

- Note that $Q(y|x) = \frac{p(x,y)}{p(x)}$. Therefore,

$$\begin{aligned}
 H(Y|X) &= -\mathbb{E}[\log_2 Q(Y|X)] = -\mathbb{E}\left[\log_2 \frac{p(X,Y)}{p(X)}\right] \\
 &= (-\mathbb{E}[\log_2 p(X,Y)]) - (-\mathbb{E}[\log_2 p(X)]) \\
 &= H(X,Y) - H(X)
 \end{aligned}$$

Example 5.5. Continue from Example 5.3. Random variables X and Y have the following joint pmf matrix \mathbf{P} :

$$\mathbf{P} = \begin{bmatrix} \frac{1}{8} & \frac{1}{16} & \frac{1}{16} & \frac{1}{4} \\ \frac{1}{16} & \frac{1}{8} & \frac{1}{16} & 0 \\ \frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \\ \frac{1}{32} & \frac{1}{32} & \frac{1}{16} & 0 \end{bmatrix}$$

Find $H(Y|X)$ and $H(X|Y)$.

$$\begin{aligned}
 H(Y|X) &= H(X,Y) - H(X) = \frac{27}{8} - \frac{7}{4} = \frac{13}{8} \\
 H(X|Y) &= H(X,Y) - H(Y) = \frac{27}{8} - 2 = \frac{11}{8}
 \end{aligned}$$

$$P = \begin{matrix} & \begin{matrix} x \backslash y & 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{bmatrix} 1/8 & 1/16 & 1/16 & 1/4 \\ 1/16 & 1/8 & 1/16 & 0 \\ 1/32 & 1/32 & 1/16 & 0 \\ 1/32 & 1/32 & 1/16 & 0 \end{bmatrix} \end{matrix} \quad \begin{matrix} p(x) \\ 1/2 \\ 1/4 \\ 1/8 \\ 1/8 \end{matrix}$$

$$Q = \begin{matrix} & \begin{matrix} x \backslash y & 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{bmatrix} 1/4 & 1/8 & 1/8 & 1/2 \\ 1/4 & 1/2 & 1/4 & 0 \\ 1/4 & 1/4 & 1/2 & 0 \\ 1/4 & 1/4 & 1/2 & 0 \end{bmatrix} \end{matrix} \quad \begin{matrix} H(Y|1) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{2}{8} \log_2 \frac{1}{8} - \frac{1}{2} \log_2 \frac{1}{2} = \frac{7}{4} \\ H(Y|2) = -\frac{3}{4} \log_2 \frac{1}{4} - \frac{1}{2} \log_2 \frac{1}{2} = 1 + \frac{1}{2} = \frac{3}{2} \\ H(Y|3) = \frac{3}{2} \\ H(Y|4) = \frac{3}{2} \end{matrix}$$

$$p(y) = \begin{matrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{matrix}$$

$$p_{x|y} = \begin{bmatrix} 1/2 & 1/4 & 1/4 & 1 \\ 1/4 & 1/2 & 1/4 & 0 \\ 1/8 & 1/8 & 1/4 & 0 \\ 1/8 & 1/8 & 1/4 & 0 \end{bmatrix}$$

$$H(X|1) = H(X|2) = H(X|3) = 2$$

$$H(X|4) = 0$$

$$H(X|Y) = \frac{1}{4} \times \frac{7}{4} + \frac{1}{4} \times \frac{7}{4} + \frac{1}{4} \times 2 + \frac{1}{4} \times 0 = \frac{11}{8}$$

Compare these with

$$P(A \cap B) = P(A) - P(A \setminus B) = P(B) - P(B \setminus A)$$

$$= P(A) + P(B) - P(A \cup B)$$

$$I(X; Y) \equiv \begin{cases} \text{The area in the middle} \\ = H(X) - H(X|Y) \\ = H(Y) - H(Y|X) \\ = H(X) + H(Y) - H(X, Y) \end{cases}$$

$$H(X) = \frac{7}{4} = \frac{14}{8}$$

$$H(Y) = 2 = \frac{16}{8}$$

$$H(X, Y) = \frac{27}{8}$$

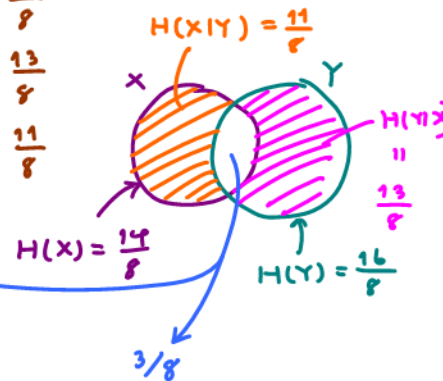
$$H(Y|X) = \frac{13}{8}$$

$$H(X|Y) = \frac{11}{8}$$

$$I(X; Y) = \frac{3}{8}$$

$$H(X, Y) = H(X) + H(Y|X)$$

$$= H(Y) + H(X|Y)$$



Definition 5.6. The **mutual information**¹⁴ $I(X; Y)$ between two random variables X and Y is defined as

$$I(X; Y) = H(X) - H(X|Y) \quad (20)$$

$$= H(Y) - H(Y|X) \quad (21)$$

$$= H(X) + H(Y) - H(X, Y) \quad (22)$$

$$= \mathbb{E} \left[\log_2 \frac{p(X, Y)}{p(X) q(Y)} \right] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x) q(y)} \quad (23)$$

$$= \mathbb{E} \left[\log_2 \frac{P_{X|Y}(X|Y)}{p(X)} \right] = \mathbb{E} \left[\log_2 \frac{Q(Y|X)}{q(Y)} \right]. \quad (24)$$

- ① Mutual information quantifies the reduction in the uncertainty of one random variable due to the knowledge of the other.

$$I(x; Y) = \underbrace{H(x)}_{\substack{\text{Amount of randomness} \\ \text{in } X}} - \underbrace{H(x|Y)}_{\substack{\text{Amount of randomness still remained in } X \\ \text{when we know } Y.}} = \text{Amount of reduction in the randomness of } X \text{ due to the knowledge of } Y.$$

- ② Mutual information is a measure of the amount of information one random variable contains about another [3, p 13].

- ③ It is natural to think of $I(X; Y)$ as a measure of distance how far X and Y are from being independent. *So, $I(X; Y)$ measures the amount of dependency btw the two RVs.*
 - Technically, it is the (Kullback-Leibler) divergence between the joint and product-of-marginal distributions.

5.7. Some important properties

- (a) $H(X, Y) = H(Y, X)$ and $I(X; Y) = I(Y; X)$. *by the symmetry in their defn.*
However, in general, $H(X|Y) \neq H(Y|X)$. *by the asymmetry in their defn.*
- (b) I and H are always ≥ 0 .
- (c) There is a one-to-one correspondence between Shannon's information measures and set theory. We may use an **information diagram**, which

¹⁴The name mutual information and the notation $I(X; Y)$ was introduced by [Fano, 1961, Ch 2].

is a variation of a Venn diagram, to represent relationship between Shannon's information measures. This is similar to the use of the Venn diagram to represent relationship between probability measures. These diagrams are shown in Figure 7.

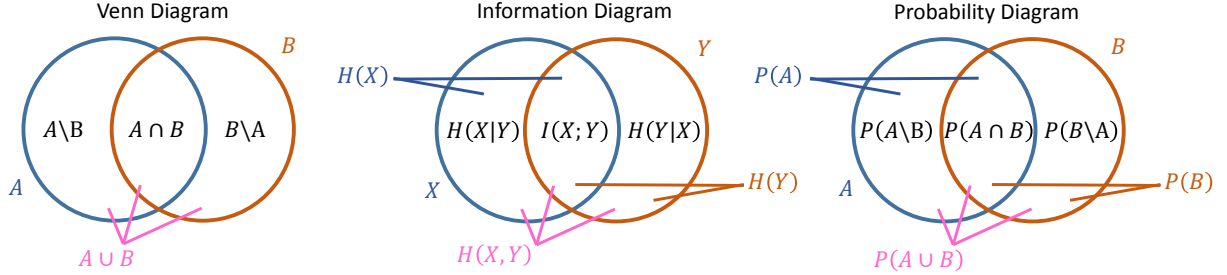


Figure 7: Venn diagram and its use to represent relationship between information measures and relationship between information measures

- Chain rule for information measures:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y).$$

(d) $I(X; Y) \geq 0$ with equality if and only if X and Y are independent.

- When this property is applied to the information diagram (or definitions (20), (21), and (22) for $I(X, Y)$), we have

- (i) $H(X|Y) \leq H(X)$,
- (ii) $H(Y|X) \leq H(Y)$,
- (iii) $H(X, Y) \leq H(X) + H(Y)$

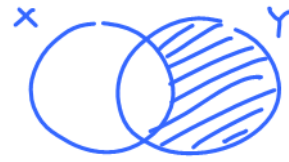
Moreover, each of the inequalities above becomes equality if and only if $X \perp\!\!\!\perp Y$.

(e) We have seen in Section 2.4 that

$$\begin{matrix} 0 \\ \text{deterministic (degenerated)} \end{matrix} \leq H(X) \leq \begin{matrix} \log_2 |\mathcal{X}| \\ \text{uniform} \end{matrix}. \quad (25)$$

Similarly,

$$\begin{matrix} 0 \\ \text{deterministic (degenerated)} \end{matrix} \leq H(Y) \leq \begin{matrix} \log_2 |\mathcal{Y}| \\ \text{uniform} \end{matrix}. \quad (26)$$



For conditional entropy, we have

$$\underset{\exists g \ Y=g(X)}{0} \leq H(Y|X) \leq \underset{X \perp\!\!\!\perp Y}{H(Y)} \quad (27)$$

and

$$\underset{\exists g \ X=g(Y)}{0} \leq H(X|Y) \leq \underset{X \perp\!\!\!\perp Y}{H(X)}. \quad (28)$$

For mutual information, we have

Think about this $\rightarrow \underset{X \perp\!\!\!\perp Y}{0} \leq I(X;Y) \leq \underset{\exists g \ X=g(Y)}{H(X)} \quad (29)$

and

$$\underset{X \perp\!\!\!\perp Y}{0} \leq I(X;Y) \leq \underset{\exists g \ Y=g(X)}{H(Y)}. \quad (30)$$

Combining 25, 26, 29, and 30, we have

$$0 \leq I(X;Y) \leq \min\{H(X), H(Y)\} \leq \min\{\log_2 |\mathcal{X}|, \log_2 |\mathcal{Y}|\} \quad (31)$$

(f) $H(X|X) = 0$ and $I(X;X) = H(X)$.

Example 5.8. Find the mutual information $I(X;Y)$ between the two random variables X and Y whose joint pmf matrix is given by $\mathbf{P} = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & 0 \end{bmatrix}$.

$$\begin{aligned} I(X;Y) &= H(X) + H(Y) - H(X,Y) = 0.1226 \\ &\quad \downarrow \quad \quad \quad \downarrow \\ &\quad 0.8113 \quad \quad \quad = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{2}{4} \log_2 \frac{1}{4} = 1.5 \\ &\quad \quad \quad \downarrow \\ &\quad \quad \quad = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8113 \end{aligned}$$

$\begin{bmatrix} \frac{1}{2} & \frac{1}{4} \\ \frac{1}{4} & 0 \end{bmatrix} \rightarrow \begin{matrix} \rightarrow 3/4 \\ \rightarrow 1/4 \end{matrix}$
 $\downarrow \quad \downarrow$
 $3/4 \quad 1/4$

Example 5.9. Find the mutual information $I(X;Y)$ between the two random variables X and Y whose $\underline{\mathbf{p}} = [\frac{1}{4}, \frac{3}{4}]$ and $\mathbf{Q} = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \\ \frac{3}{4} & \frac{1}{4} \end{bmatrix}$.

$$I(X;Y) = 0.1432$$

$$\begin{bmatrix} \frac{1}{4} & \frac{3}{4} \\ \frac{3}{4} & \frac{1}{4} \end{bmatrix} \xrightarrow{\times 1/4} \begin{bmatrix} \frac{1}{16} & \frac{3}{16} \\ \frac{9}{16} & \frac{3}{16} \end{bmatrix} = \mathbf{P}$$

$\times 3/4$

$$P = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \end{bmatrix}$$

$$Q = \begin{array}{c|cc} x \backslash y & 0 & 1 \\ \hline 0 & \frac{1}{4} & \frac{3}{4} \\ 1 & \frac{3}{4} & \frac{1}{4} \end{array}$$

Method 1 : $I(X;Y) = H(X) + H(Y) - H(X,Y)$

First, convert the given information into the joint pmf matrix.

$$P = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \end{bmatrix} \rightarrow H(X) = H\left(\begin{bmatrix} \frac{1}{4} & \frac{3}{4} \end{bmatrix}\right) = 0.8113$$

$$Q = \begin{array}{c|cc} x \backslash y & 0 & 1 \\ \hline 0 & \frac{1}{4} & \frac{3}{4} \\ 1 & \frac{3}{4} & \frac{1}{4} \end{array} \xrightarrow{\begin{array}{l} \times \frac{1}{4} \\ \times \frac{3}{4} \end{array}} \begin{array}{c|cc} x \backslash y & 0 & 1 \\ \hline 0 & \frac{1}{16} & \frac{9}{16} \\ 1 & \frac{9}{16} & \frac{3}{16} \end{array} = P \left\{ \begin{array}{l} \rightarrow H(X,Y) = H\left(\begin{bmatrix} \frac{1}{16} & \frac{9}{16} & \frac{9}{16} & \frac{3}{16} \end{bmatrix}\right) \\ \quad \quad \quad = 1.6226 \\ \downarrow \quad \downarrow \text{ (sum along each column)} \\ \frac{10}{16} \quad \frac{6}{16} \rightarrow H(Y) = H\left(\begin{bmatrix} \frac{5}{8} & \frac{3}{8} \end{bmatrix}\right) = 0.9544 \end{array} \right.$$

$$I(X;Y) = H(X) + H(Y) - H(X,Y) \approx 0.1432.$$

Method 2 : $I(X;Y) = H(Y) - H(Y|X)$

$$H(Y|X) = \sum_x p(x) \underbrace{H(Y|x)}_{0.8113} = 0.8113 = \sum_x \overbrace{p(x)}^1 = 0.8113$$

0.8113 ← Each row of Q contains $\frac{1}{4}, \frac{3}{4}$. So, $H(Y|x) = 0.8113$ for any x . (for any row.)

Same as Ex. 5.8

$$I(X;Y) = \underbrace{H(Y)}_{0.9544} - H(Y|X) = 0.1432$$

To find $H(Y)$, we need $q(y)$

$$\begin{aligned} \underline{q} &= P \cdot Q = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \end{bmatrix} \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \\ \frac{3}{4} & \frac{1}{4} \end{bmatrix} \\ &= \begin{bmatrix} \frac{10}{16} & \frac{6}{16} \end{bmatrix} = \begin{bmatrix} \frac{5}{8} & \frac{3}{8} \end{bmatrix} \end{aligned}$$